

Informative Data Projections: A Framework and Two Examples

Tijl De Bie^{1,2}, Jefrey Lijffijt^{1,2}, Raúl Santos-Rodriguez², and Bo Kang¹ *

1- Data Science Lab - Ghent University
Technicum, Sint-Pietersnieuwstraat 41, 9000 Gent - Belgium

2- Dept. of Engineering Mathematics - University of Bristol
MVB Woodland Road, BS8 1UB, Bristol - United Kingdom

Abstract. Methods for Projection Pursuit aim to facilitate the visual exploration of high-dimensional data by identifying interesting low-dimensional projections. A major challenge is the design of a suitable quality metric of projections, commonly referred to as the *projection index*, to be maximized by the Projection Pursuit algorithm. In this paper, we introduce a new information-theoretic strategy for tackling this problem, based on quantifying the amount of information the projection conveys to a user given their *prior beliefs* about the data. The resulting projection index is a subjective quantity, explicitly dependent on the intended user. As a useful illustration, we developed this idea for two particular kinds of prior beliefs. The first kind leads to PCA (Principal Component Analysis), shining new light on when PCA is (not) appropriate. The second kind leads to a novel projection index, the maximization of which can be regarded as a robust variant of PCA. We show how this projection index, though non-convex, can be effectively maximized using a modified power method as well as using a semidefinite programming relaxation. The usefulness of this new projection index is demonstrated in comparative empirical experiments against PCA and a popular Projection Pursuit method.

1 Introduction

The analysis of high-dimensional data often starts with dimensionality reduction, to facilitate initial visual exploration by a human user. Most analysts will instinctively do this by Principal Component Analysis (PCA) [1]: it is widely available, computationally efficient, easy to interpret, and in the common situation where the data lies close to a low-dimensional subspace, PCA is effective in retrieving it. However, in user interactions with their PRIM-9 system for interactive data exploration [2], it was observed that the human operators tended to prefer projections that reveal *some form of structure*, rather than projections of *high variance* as preferred by PCA. Later [3] provided theoretical arguments for why projections in which the data are Normally distributed are *least* interesting, as they essentially reveal no structure in the data.¹

*This work was supported by the European Union through the ERC Consolidator Grant FORSIED (project reference 615517).

¹This means that Independent Component Analysis (ICA) and Projection Pursuit (PP) are largely equivalent: appropriate PP methods can be and are being used to do ICA [4, 5, 6].

Quantifying the precise extent to which a projection *is* interesting, however, is riddled with conceptual and practical difficulties. In fact, it seemed obvious to the early PP research protagonists that a universally useful *projection index* that formalizes the interestingness of a projection cannot exist (see e.g. [3]). The answer, therefore, was the introduction of lots of different projection indices. Most of these aim to quantify the extent to which the distribution of the projected data departs from the Normal distribution, and all strike a different balance between practical usefulness, computational complexity, and robustness against outliers (e.g. [2, 3, 4, 7] and references therein). Indeed, due to the elusive nature of the core question of what makes a projection interesting to a given user, the focus shifted towards secondary questions around robustness aspects and computational properties of the projection indices.

Contributions in this paper Here our aim is to return the focus to the user once again, and directly ask the question of how interesting a given data projection is *to a particular user*. Our work presents the first generic design strategy for projection indices that explicitly depend on the intended user.

In Section 2 we introduce a strategy for quantifying the interestingness of a projection as its information content against the background of a probability distribution representing the user’s beliefs about the data. In Sections 3 and 4 we then apply this strategy for two particular types of prior beliefs, leading to a novel interpretation for PCA in the first case, and a novel projection index in the second, the optimization of which represents a robust variant of PCA. Although this latter projection index is non-convex, we introduce two algorithms for effectively optimizing it. We end by empirically illustrating the benefits of this robust PCA variant as compared to standard PCA and FastICA, a popular PP method that is also used for ICA [6].

2 The subjective information content of a data projection – general outline

The proposed strategy for quantifying the subjective information content of a data projection closely follows the generic approach introduced by [8, 9], where the choices made are extensively motivated. Here we will merely provide some intuition underlying the approach.

The strategy from [8] for quantifying information content rests on the availability of a representation of the user’s prior belief state in the form of a probability density $p_{\mathbf{X}}$ over the set of possible values for the data \mathbf{X} – *in casu* over the set $\mathbb{R}^{n \times d}$. Given this so-called *background distribution*, one can then compute the marginal probability density function of a data projection $\mathbf{p}_{\mathbf{w}} = \mathbf{X}\mathbf{w}$ defined by the weight vector $\mathbf{w} \in \mathbb{R}^d$. We will denote this marginal probability density function as $p_{\mathbf{X}\mathbf{w}}$.

We call a *projection pattern* a statement of the form $\mathbf{p}_{\mathbf{w}} \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1})$, specifying that the value $\mathbf{p}_{\mathbf{w}}$ of the projected data lies within an interval of width Δ around $\hat{\mathbf{X}}\mathbf{w}$ (with $\hat{\mathbf{X}} \in \mathbb{R}^{n \times d}$ the empirical data). This is what is conveyed to a

user through a scatter plot of the data projections $\hat{\mathbf{X}}\mathbf{w}$, with plotting resolution Δ . Clearly, the smaller the probability $\text{Prob}_{\mathbf{p}_w \sim p_{\mathbf{X}w}}(\mathbf{p}_w \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1}])$, the more the surprising and hence informative this pattern would be *to that particular user*. This is argued more formally by [8], where the negative logarithm of this probability is shown to be a good measure of Subjective Information Content (SIC). We denote this as follows:

$$\text{SIC}(\hat{\mathbf{X}}\mathbf{w}) = -\log\left(\text{Prob}_{\mathbf{p}_w \sim p_{\mathbf{X}w}}(\mathbf{p}_w \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1}])\right)$$

This is what we propose as a generic projection index, quantifying the interestingness of a projection.

An important question is how $p_{\mathbf{X}}$ and hence the marginals $p_{\mathbf{X}w}$ can be obtained, without overburdening the user. In [8] it is suggested that the user is often capable of specifying aspects of their belief state as constraints on expected values of specified statistics of the data. He argued that the *Maximum Entropy* (MaxEnt) distribution subject to these constraints is an attractive choice, given its unbiasedness, its robustness, and in being an *exponential family* model [10], the inference of which is well understood and often computationally tractable.

In Sections 3 and 4, we will develop this strategy for two different assumptions regarding the prior beliefs of the user, illustrating how it can lead to new algorithms, as well as to new insights into existing ones. Throughout this paper, we assume the data has been centered (i.e. has zero mean).

3 PCA: an information theoretic interpretation

Here we show how standard PCA can be derived using this generic strategy for designing projection indices. This exercise will also present new insight in what cases PCA is an effective PP approach.

3.1 The prior beliefs

A user not expecting any outliers can be assumed capable of expressing an expectation about the value of the average two-norm squared of the data points:

$$\mathbb{E}_{\mathbf{X} \sim p_{\mathbf{X}}} \left\{ \frac{1}{n} \sum_i^n \mathbf{x}_i' \mathbf{x}_i \right\} = \sigma^2. \quad (1)$$

I.e., the user has a specific expectation about the average squared norm of the data points.

The MaxEnt distribution subject to this constraint is well known and equal to a product distribution of multivariate Normal distributions $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, with one factor for each of the data points \mathbf{x}_i . More formally, the density function $p_{\mathbf{X}}$ for the dataset representing the background distribution is:

$$p_{\mathbf{X}}(\mathbf{X}) = \prod_i p_{\mathbf{x}}(\mathbf{x}_i), \text{ where } p_{\mathbf{x}}, \text{ defined as } p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}^d} \exp\left(-\frac{\mathbf{x}'\mathbf{x}}{2\sigma^2}\right), \quad (2)$$

is the probability density function for each of the individual data points in the dataset.²

3.2 The subjective information content

Given a Normal random vector $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, a projection onto a weight vector \mathbf{w} with $\mathbf{w}'\mathbf{w} = 1$ is also Normal: $\mathbf{x}'\mathbf{w} \sim \mathcal{N}(0, \sigma)$. Thus, given the independence of the data points under the background distribution, the marginal probability density function $p_{\mathbf{X}\mathbf{w}}$ for the projection $\mathbf{p}_{\mathbf{w}} = \mathbf{X}\mathbf{w}$ of a dataset \mathbf{X} sampled from the background distribution is given by:

$$p_{\mathbf{X}\mathbf{w}}(\mathbf{p}_{\mathbf{w}}) = \frac{1}{\sqrt{2\pi\sigma^2}^n} \exp\left(-\frac{\mathbf{p}_{\mathbf{w}}'\mathbf{p}_{\mathbf{w}}}{2\sigma^2}\right).$$

We can thus compute the SIC of a projection pattern $\mathbf{p}_{\mathbf{w}} \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1}]$ as minus the logarithm of its probability under this marginal density function $p_{\mathbf{X}\mathbf{w}}$. Noting that for small enough Δ , $\text{Prob}_{\mathbf{p}_{\mathbf{w}} \sim p_{\mathbf{X}\mathbf{w}}}(\mathbf{p}_{\mathbf{w}} \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1}]) \approx \Delta^n \cdot p_{\mathbf{X}\mathbf{w}}(\hat{\mathbf{X}}\mathbf{w})$, this leads to:

$$\begin{aligned} \text{SIC}(\hat{\mathbf{X}}\mathbf{w}) &= -\log(p_{\mathbf{X}\mathbf{w}}(\hat{\mathbf{X}}\mathbf{w})) - n \log(\Delta) \\ &= \frac{n}{2} \log(2\pi\sigma^2) - n \log(\Delta) + \frac{1}{2\sigma^2} \mathbf{w}'\hat{\mathbf{X}}'\hat{\mathbf{X}}\mathbf{w}. \end{aligned} \quad (3)$$

It is trivial to generalize this toward r -dimensional projections $\mathbf{P}_{\mathbf{W}_r} = \mathbf{X}\mathbf{W}_r$ of the dataset, defined by an orthogonal matrix $\mathbf{W}_r \in \mathbb{R}^{d \times r}$ with $\mathbf{W}_r'\mathbf{W}_r = \mathbf{I}$. With $\Delta' = (\Delta_1 \ \Delta_2 \ \dots \ \Delta_r)$ a vector containing the resolutions for each of the projections, then the information content of the pattern $\mathbf{P}_{\mathbf{W}_r} \in [\hat{\mathbf{X}}\mathbf{W}_r, \hat{\mathbf{X}}\mathbf{W}_r + \mathbf{1}\Delta']$ specifying r projections simultaneously, is given by:

$$\text{SIC}(\hat{\mathbf{X}}\mathbf{W}_r) = \frac{nr}{2} \log(2\pi\sigma^2) - n \sum_{i=1}^r \log(\Delta_i) + \frac{1}{2\sigma^2} \text{Tr}[\mathbf{W}_r'\hat{\mathbf{X}}'\hat{\mathbf{X}}\mathbf{W}_r]. \quad (4)$$

3.3 Finding the most informative projections

For fixed Δ , maximizing the SIC from Eq. (3) is done by solving:

$$\max_{\mathbf{w}} \mathbf{w}'\hat{\mathbf{X}}'\hat{\mathbf{X}}\mathbf{w}, \text{ s. t. } \mathbf{w}'\mathbf{w} = 1,$$

equivalent to the optimization problem to be solved for finding the first principal component in classical PCA. Similarly, optimizing Eq. (4) is equivalent to finding the r dominant PCA components.

Remark 1. *The assumption that Δ is constant is not always warranted. When making a scatter plot, it is often desirable to stretch the axes in order to fill*

²Note that it is not our intention to argue in favor of this; our intent is merely to investigate what is a suitable projection index if this is an accurate representation of the prior belief state.

available space. Then $\Delta \propto \max(\hat{\mathbf{X}}\mathbf{w}) - \min(\hat{\mathbf{X}}\mathbf{w})$, such that finding the most informative projection amounts to solving:

$$\max_{\mathbf{w}} \mathbf{w}'\hat{\mathbf{X}}'\hat{\mathbf{X}}\mathbf{w} - 2\sigma^2 n \log \left(\max(\hat{\mathbf{X}}\mathbf{w}) - \min(\hat{\mathbf{X}}\mathbf{w}) \right), \text{ s. t. } \mathbf{w}'\mathbf{w} = 1,$$

This is another explanation of why simple variance maximization done by PCA is rarely a good approach for exploratory data analysis.

From Eq. (4) the most informative r -dimensional projection of the data are found by solving:

$$\max_{\mathbf{W}_r \in \mathbb{R}^{d \times r}} \text{Tr} \left[\mathbf{W}_r' \hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{W}_r \right], \text{ s. t. } \mathbf{W}_r' \mathbf{W}_r = \mathbf{I},$$

again equivalent to the optimization problem for finding the r dominant principal components.

Remark 2. *The exact value of σ will have no effect on the relative information content of different possible projections, hence the user does not need to provide this value.*

4 t-PCA: for users expecting a heavy tailed distribution

The previous section elucidates the assumptions on the user (prior belief on average squared norm of the data points) and visualization approach (constant resolution) for PCA to be optimal. In the present section we will develop an alternative for PCA when the assumption on the user's prior beliefs is altered, to be more accommodating for outliers.

4.1 The prior beliefs

As the user's prior belief about the data, we here propose the following:

$$\mathbb{E}_{\mathbf{X} \sim p_{\mathbf{X}}} \left\{ \frac{1}{n} \sum_i \log \left(1 + \frac{1}{\rho} \mathbf{x}_i' \mathbf{x}_i \right) \right\} = c.$$

Thus, rather than specifying an expectation on the spread of the data, for small values of ρ the user specifies an expectation on the *order of magnitude* of the spread of the data. When the user expects outliers to be present, they may feel able to specify an expectation on the average *order of magnitude* of the 2-norms of the data points, rather than on the average of their 2-norms themselves.

For notational convenience, let us introduce the function $\kappa(\nu) = \psi\left(\frac{\nu+d}{2}\right) - \psi\left(\frac{\nu}{2}\right)$, where ψ represents the digamma function. In the sequel the value of $\kappa^{-1}(c)$ will need to be used, denoted as ν for brevity. Then, the background distribution can be derived by relying on [11], where it is shown that the MaxEnt

distribution subject to the specified prior information is the product of independent multivariate standard t -distributions with density function $p_{\mathbf{x}}$ defined as:

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{\Gamma\left(\frac{\nu+d}{2}\right)}{\sqrt{(\pi\rho)^d}\Gamma\left(\frac{\nu}{2}\right)} \cdot \frac{1}{\left(1 + \frac{1}{\rho}\mathbf{x}'\mathbf{x}\right)^{\frac{\nu+d}{2}}},$$

with a factor in this product distribution for each data point. Here Γ represents the gamma function.

Note that for $\rho, \nu \rightarrow \infty$, $\frac{\rho}{\nu} \rightarrow \sigma^2$ this density function tends to the multivariate Normal density function with mean $\mathbf{0}$ and covariance $\sigma^2\mathbf{I}$. For $\rho = \nu = 1$ it is a multivariate standard Cauchy distribution, which is so heavy-tailed that its mean is undefined and its second moment is infinitely large. Thus, this type of prior belief can clearly model the expectation of outliers to varying degrees.

4.2 The subjective information content

The marginals of a t -distribution with given correlation matrix are again a t -distribution with the same number of degrees of freedom, obtained by simply selecting the relevant part of the correlation matrix [12, 13]. This means that the marginal density function for the data projections $\mathbf{p}_{\mathbf{w}} = \mathbf{X}\mathbf{w}$ onto a vector \mathbf{w} with $\mathbf{w}'\mathbf{w} = 1$ (and $p_{\mathbf{w},i} \triangleq \mathbf{x}'_i\mathbf{w}$) is:

$$p_{\mathbf{X}\mathbf{w}}(\mathbf{p}_{\mathbf{w}}) = \prod_i p_{\mathbf{x}'\mathbf{w}}(p_{\mathbf{w},i}), \text{ where } p_{\mathbf{x}'\mathbf{w}}(p_{\mathbf{w},i}) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\rho}\Gamma\left(\frac{\nu}{2}\right)} \cdot \frac{1}{\left(1 + \frac{1}{\rho}p_{\mathbf{w},i}^2\right)^{\frac{\nu+1}{2}}}.$$

Thus the SIC of the projection pattern $\mathbf{p}_{\mathbf{w}} \in [\hat{\mathbf{X}}\mathbf{w}, \hat{\mathbf{X}}\mathbf{w} + \Delta\mathbf{1}]$ is:

$$\text{SIC}\left(\hat{\mathbf{X}}\mathbf{w}\right) = \frac{\nu+1}{2} \sum_{i=1}^n \log\left(1 + \frac{1}{\rho}(\hat{\mathbf{x}}'_i\mathbf{w})^2\right) - n \log(\Delta) + \text{a constant}. \quad (5)$$

This derivation can be generalized towards the information content of an r -dimensional projection onto the columns of an orthogonal matrix $\mathbf{W}_r \in \mathbb{R}^{d \times r}$ (i.e. $\mathbf{W}_r'\mathbf{W}_r = \mathbf{I}$). Without proof, we state the information content of the pattern $\mathbf{P}_{\mathbf{W}_r} \in [\hat{\mathbf{X}}\mathbf{W}_r, \hat{\mathbf{X}}\mathbf{W}_r + \mathbf{1}\Delta']$:

$$\text{SIC}\left(\hat{\mathbf{X}}\mathbf{W}_r\right) = \frac{\nu+r}{2} \sum_{i=1}^n \log\left(1 + \frac{1}{\rho}\hat{\mathbf{x}}'_i\mathbf{W}_r\mathbf{W}_r'\hat{\mathbf{x}}_i\right) - n \sum_{i=1}^r \log(\Delta_i) + \text{a constant}. \quad (6)$$

4.3 Finding the most informative projections: an analysis and two algorithms

As in the derivation of PCA, we will assume that Δ is constant. For ease of exposition, we will focus on the SIC of a single projection as given by Eq. (5).

Taking into account that $\mathbf{w}'\mathbf{w} = 1$, and ignoring some constant factors and terms, maximising the SIC is thus equivalent to solving the following problem:

$$\max_{\mathbf{w}} \sum_{i=1}^n \log(\rho + (\hat{\mathbf{x}}'_i \mathbf{w})^2), \text{ s. t. } \mathbf{w}'\mathbf{w} = 1. \quad (7)$$

Clearly, the larger $\mathbf{w}'\mathbf{w}$, the larger the objective, so the constraint can be relaxed to $\mathbf{w}'\mathbf{w} \leq 1$.

Remark 3. *Given the reliance of this approach on the multivariate t -distribution as a background distribution, we will refer to this approach as t -PCA.*

Remark 4. *Just like in PCA where the value of σ has no effect on which pattern is most interesting, here the value of ν and thus of c is absent from the final optimization problem, and thus it has no effect on which projection is the most interesting one. (Though σ and c do affect the value of the interestingness.) This significantly reduces the demands on the user in specifying their prior beliefs.*

Remark 5. *By varying ρ , t -PCA interpolates between maximizing the arithmetic mean, like PCA does, and maximizing the geometric mean of the squares of the data projections, which is more robust against outliers. Indeed, for $\rho = 0$, the objective function is monotonically related to the geometric mean of the squares of the data projections $(\hat{\mathbf{x}}_i \mathbf{w})^2$:*

$$\exp \left[\frac{1}{n} \sum_{i=1}^n \log(\hat{\mathbf{x}}'_i \mathbf{w})^2 \right] = \left(\prod_{i=1}^n (\hat{\mathbf{x}}'_i \mathbf{w})^2 \right)^{\frac{1}{n}}.$$

On the other hand, for $\rho \rightarrow \infty$, the objective function is monotonically related to the arithmetic mean, and thus becomes equivalent to the PCA objective function:

$$\lim_{\rho \rightarrow \infty} \frac{\rho}{n} \sum_{i=1}^n \log(\rho + (\hat{\mathbf{x}}'_i \mathbf{w})^2) - \rho \log(\rho) = \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{x}}'_i \mathbf{w})^2.$$

4.3.1 The complexity of the optimization problem

To get some insight into the computational complexity of problem 7, let us consider the special case of $\rho = 0$. The constraint $\mathbf{w}'\mathbf{w} \leq 1$ is convex, and the objective is concave as long as $\mathbf{w} \in \mathcal{W}_{\mathbf{s}} \triangleq \{\mathbf{w} | \text{sign}(\hat{\mathbf{X}}\mathbf{w}) = \mathbf{s}\}$ for some fixed sign vector \mathbf{s} . Indeed, for $\rho = 0$ the objective function can be rewritten as $\sum_{i=1}^n \log((\hat{\mathbf{x}}'_i \mathbf{w})^2) = \sum_{i=1}^n \log \det \begin{pmatrix} s_i \hat{\mathbf{x}}'_i \mathbf{w} & 0 \\ 0 & s_i \hat{\mathbf{x}}'_i \mathbf{w} \end{pmatrix}$, which is the sum of n (concave) log determinant functions of linear matrix functions of the parameters \mathbf{w} .

This seems to suggest a possible solution strategy, at least for the case $\rho = 0$: enumerate all possible sign vectors \mathbf{s} for the dataset $\hat{\mathbf{X}}$, find a weight vector \mathbf{w} for each of these, and locally optimize it using a convex optimization problem.

However, according to Cover's Function-Counting Theorem [14], the number of homogeneously linearly separable dichotomies of n points in d -dimensional Euclidean space is $2 \sum_{k=0}^{d-1} \binom{n-1}{k} = O((n-1)^{d-1})$, which is clearly impractical.

Since even for the special case of $\rho = 0$ it is impractical to maximize the information content exactly, we developed two approximation algorithms: the first one a modification of the power method for solving eigenvalue problems, and the second one a convex relaxation to a log-determinant optimization problem.

4.3.2 A modified power method

The stationarity condition of the Karush-Kuhn-Tucker (KKT) optimality conditions for Eq. (7) is:

$$\left(\sum_{i=1}^n \frac{\hat{\mathbf{x}}_i \hat{\mathbf{x}}_i'}{\rho + (\hat{\mathbf{x}}_i' \mathbf{w})^2} \right) \mathbf{w} = \lambda \mathbf{w},$$

with $\lambda \geq 0$ a KKT multiplier corresponding to the constraint $\mathbf{w}'\mathbf{w} \leq 1$. Note that the matrix on the left hand side is essentially a weighted empirical covariance matrix for the data, where points contribute more if they have a smaller value for $(\hat{\mathbf{x}}_i' \mathbf{w})^2$: the weight for $\hat{\mathbf{x}}_i \hat{\mathbf{x}}_i'$ is $\frac{1}{\rho + (\hat{\mathbf{x}}_i' \mathbf{w})^2}$.

As pointed out above, this optimisation problem is not convex, and also the optimality conditions do not admit a closed form solution in terms of e.g. an eigenvalue problem. Instead we investigated the use of a simple gradient descent method with after each gradient step a projection onto the feasible set $\mathbf{w}'\mathbf{w} = 1$. This can be formulated as a modified power method [15]:

1. Start with an initial value for $\mathbf{w}^{(0)}$, normalized to unit norm.
2. Iterate from $k = 1$ until convergence or maximum number of iterations reached:
 - (a) $\mathbf{v}^{(k)} = \left(\mathbf{I} + \alpha \sum_{i=1}^n \frac{\hat{\mathbf{x}}_i \hat{\mathbf{x}}_i'}{\rho + (\hat{\mathbf{x}}_i' \mathbf{w}^{(k-1)})^2} \right) \mathbf{w}^{(k-1)}$.
 - (b) $\mathbf{w}^{(k)} = \frac{\mathbf{v}^{(k)}}{\|\mathbf{v}^{(k)}\|}$.

Here, α is a step-size parameter that controls the speed of convergence. Clearly this algorithm is not guaranteed to converge, but for sufficiently small α it does converge to a local optimum in practice.³

For $\mathbf{w}^{(0)}$, we use the dominant eigenvector of $\left(\sum_{i=1}^n \frac{\hat{\mathbf{x}}_i \hat{\mathbf{x}}_i'}{\rho + \hat{\mathbf{x}}_i' \hat{\mathbf{x}}_i} \right)$, which amounts to maximizing an approximation of the SIC obtained by approximating $(\hat{\mathbf{x}}_i' \mathbf{w})^2$ with $\hat{\mathbf{x}}_i' \hat{\mathbf{x}}_i$.

As a first approximation, we search for good subsequent projections simply by iteratively deflating the data (i.e. projecting it onto the orthogonal complement of the previously found projections).

³A detailed convergence analysis is left as further work.

4.3.3 A convex relaxation

The problem can also be relaxed to a semidefinite log determinant optimization problem. We will do this more generally for finding the most informative r -dimensional projection for general r , i.e. of the optimization problem maximizing Eq. (6) subject to $\mathbf{W}_r' \mathbf{W}_r = \mathbf{I}$. After removing irrelevant constant terms and factors, the optimization problem we need to solve is equivalent with:

$$\max_{\mathbf{W}_r \in \mathbb{R}^{d \times r}} \sum_{i=1}^n \log(\rho + \mathbf{x}_i' \mathbf{W}_r \mathbf{W}_r' \mathbf{x}_i), \text{ s. t. } \mathbf{W}_r' \mathbf{W}_r = \mathbf{I}. \quad (8)$$

We claim that this problem can be rewritten in terms of a new variable $\mathbf{M} \triangleq \mathbf{W}_r \mathbf{W}_r'$ as follows:

$$\max_{\mathbf{M} \in \mathbb{R}^{d \times d}} \sum_{i=1}^n \log(\rho + \mathbf{x}_i' \mathbf{M} \mathbf{x}_i), \text{ s. t. } \begin{cases} \text{Tr}[\mathbf{M}] = r, \\ \mathbf{M} \succeq \mathbf{0}, \\ \mathbf{I} - \mathbf{M} \succeq \mathbf{0}, \\ \text{Rank}(\mathbf{M}) = r. \end{cases} \quad (9)$$

Theorem 1. *Problems 8 and 9 are equivalent.*

Proof. Given $\mathbf{M} = \mathbf{W}_r \mathbf{W}_r'$, the objective functions are clearly equivalent. Thus all we need to show is that the feasible sets are identical as well. It is easy to verify that the constraints in problem 8 imply the constraints in problem 9. Also the converse is true. Indeed, the rank constraint forces all but r eigenvalues to be 0. The constraints $\mathbf{I} - \mathbf{M} \succeq \mathbf{0}$ and $\mathbf{M} \succeq \mathbf{0}$ ensure all eigenvalues lie between 0 and 1. And the trace constraint ensures all eigenvalues add up to r . Hence, \mathbf{M} must have r eigenvalues equal to 1, and $d - r$ equal to 0. Thus, we can write the eigenvalue decomposition of \mathbf{M} as: $\mathbf{M} = \mathbf{W}_r \mathbf{I} \mathbf{W}_r'$, with \mathbf{W}_r an orthogonal matrix as required by the constraint in the previous formulation. Thus, the constraint set of problem 9 also implies the constraint of problem 8. \square

The only non-convex constraint in problem formulation 9 is the rank constraint. We suggest to drop that constraint to relax this problem. Note that for $r = 1$, the constraint $\mathbf{I} - \mathbf{M} \succeq \mathbf{0}$ is redundant given $\text{Tr}[\mathbf{M}] = 1$ and $\mathbf{M} \succeq \mathbf{0}$, and can therefore also be dropped. To obtain an estimate for the unrelaxed weight matrix \mathbf{W}_r , we suggest to use the r dominant eigenvectors of \mathbf{M} .

5 Empirical evaluation

In Section 5.1 we evaluate the behaviour of t-PCA in comparison with PCA in a controlled setting on synthetic data. Then, in Section 5.2, we demonstrate the practical usefulness of t-PCA on some real-life datasets, and compare its results with PCA as well as a popular projection pursuit method (FastICA).

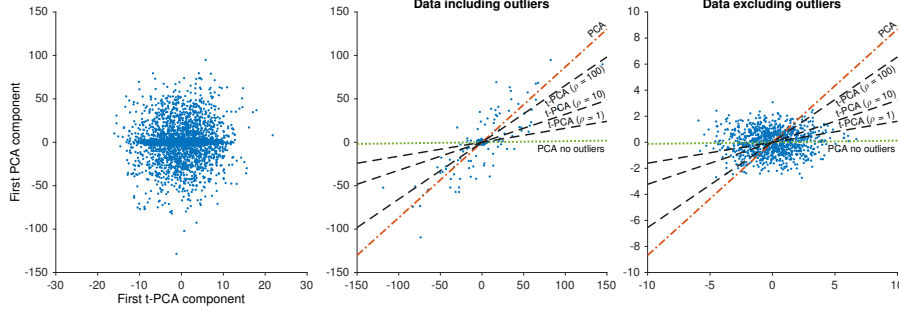


Fig. 1: Left: comparison of the dominant PCA projection (vertical axis) with the most informative t-PCA projection with $\rho = 0$ (horizontal axis). t-PCA shows more detail in the central less spread out point cloud that contains most data points. Middle and right: a scatter plot of all data points in the original space, including (middle) and excluding (right) outliers, with weight vectors of PCA (dot-dashed red line), as well as t-PCA computed with the modified power method with $\rho = 1, 10, 100$ (dashed black lines) and PCA fitted on data excluding the 100 outliers (dotted green line).

5.1 Experiments on synthetic data

Comparing PCA with t-PCA We generated a synthetic dataset consisting of two populations: a population with a small spread and 8000 data points, and a population with a large spread and 2000 data points. More specifically, both populations were sampled from a 100-dimensional multivariate Normal distribution with diagonal covariance. For the large population the variances were sampled from a χ^2 distribution with 1 degree of freedom; for the small population the same process was used, but the covariance matrix was then multiplied by 100.

As shown in Fig. 1 (left figure), due to the sensitivity of PCA to outliers, the dominant PCA direction is determined almost exclusively by the small population with large spread. t-PCA (here with the power method), however, offers an insight into the large population with lower spread as well.

The information content of the solution found by the SDP relaxation and modified power method are 1.43×10^4 and 1.44×10^4 respectively, while that of PCA is much worse at 3.18×10^3 . Note that the SDP relaxation took around 3 hours, compared to 30 seconds for the modified power method, on this $10,000 \times 100$ dataset. Usefully, the relaxation also provides us with an upper bound, namely 1.681×10^4 . Thus, both the SDP relaxation and the modified power method are close to optimal.

The effect of ρ To illustrate the robustness of t-PCA, consider a dataset consisting of two populations with different covariance structures: 1000 data points

sampled from the Normal distribution $\mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}\right)$, and 100 ‘outliers’ from a Normal distribution $\mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 16 & 12 \\ 12 & 13 \end{pmatrix}\right)$.

The weight vector resulting from PCA is shown with a dashed red line in Fig. 1 (right two figures). The full black lines show the weight vectors retrieved by t-PCA, with values for ρ equal to 1, 10, and 100. The largest value of these resulted in the line closest to the PCA result. The green dotted line shows the weight vector that would have been found using PCA had there been no outliers at all (i.e. computed just on the first 1000 data points). The middle plot in Fig. 1 demonstrates that the PCA result is determined primarily by the outliers. The right plot shows the same resulting weight vectors on top of a scatter plot of excluding the 100 outliers, showing that t-PCA is hardly affected by the outliers.

5.2 Experiments on real-life data

We used two realistic datasets: the Shuttle Dataset⁴ (58000 datapoints and 9 numerical dimensions) available from the UCI repository, and a reduced version⁵ of the 20 NewsGroups dataset (16242 datapoints and 100 dimensions). Both these datasets exhibit some complex structure, and the former in particular has a highly imbalanced cluster structure (class 1 contains 80% of all data points).

The algorithms evaluated include PCA, t-PCA with ρ equal to 10^{-5} multiplied by a measure of the scale of the data equal to the square root of the average squared norm of all data points, and a popular PP method often used for ICA, known as FastICA, with default parameters.

The results are shown in Fig. 2. The four top-level newsgroup classes in the reduced 20 NewsGroups dataset, and the seven classes in the Shuttle dataset, are shown in different colours. In all cases, the t-PCA version appears to reveal a more interesting structure in the data than either PCA or FastICA do. Remarkably, for the 20 NewsGroups dataset, the t-PCA weight vectors are close to sparse: more than 97% of the total variance of the weight vectors are due to just three, four of the 100 dimensions (i.e. words), respectively: ‘email’, ‘help’, and ‘problem’ for the first projection and ‘case’, ‘fact’, ‘god’, and ‘question’ for the second. The FastICA weight vectors, in contrast, have almost all weight on a single dimension, explaining that all datapoints are projected onto one of just three points when projecting onto the two top ICA components.

6 Conclusions

PCA is often notoriously inappropriate for dimensionality reduction, e.g. in the presence of outliers. To address this the Projection Pursuit literature has introduced numerous *projection indices* that quantify the interestingness of a projection in different ways. More recently, various authors also proposed principled *robust* versions of PCA as an alternative, e.g. [16, 17]. Yet, while these

⁴[https://archive.ics.uci.edu/ml/datasets/Statlog+\(Shuttle\)](https://archive.ics.uci.edu/ml/datasets/Statlog+(Shuttle))

⁵<http://cs.nyu.edu/roweis/data.html>

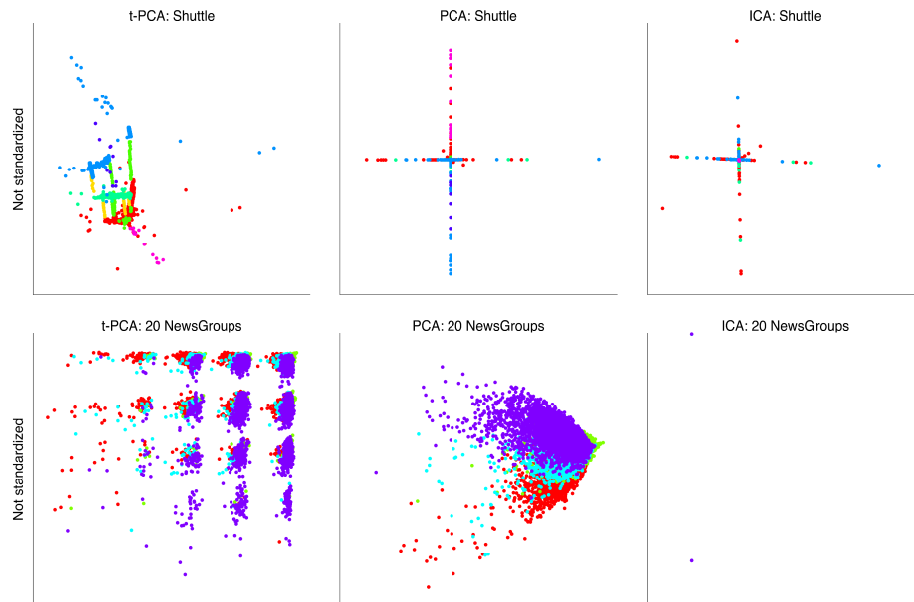


Fig. 2: The top 2 projections found by t-PCA (left), PCA (middle), and FastICA (right). Top row: Shuttle; bottom row: 20 NewsGroups.

lines of work are useful when the assumptions made are valid, they do not fundamentally address how interesting a data projection is to a user. We presented a new approach to this elusive problem, explicitly recognizing the subjective nature of the notion ‘interestingness’.

Avenues for further work include alternative prior beliefs and data types, e.g. the case where the data is not real-valued but positive integer-valued or where the user assumes dependencies between the data points (when they are vectors in a time series, geographical locations, people in a social network, etc.). More immediately, the computational properties of the t-PCA optimization problem and its convex relaxation are worth investigating. Finally, an open question is to what extent the proposed strategy can be applied to non-linear dimensionality reduction as well.

References

- [1] Ian Jolliffe. *Principal component analysis*. Wiley Online Library, 2002.
- [2] Jerome H Friedman and John W Tukey. A projection pursuit algorithm for exploratory data analysis. *IEEE Tr. Comp.*, 100(23), 1974.
- [3] Peter J Huber. Projection pursuit. *Ann. Stat.*, 13(2):435–475, 1985.
- [4] Aapo Hyvärinen. New approximations of differential entropy for independent component analysis and projection pursuit. In *Proc. of NIPS*, volume 10, page 273. MIT Press, 1998.

- [5] Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent component analysis*, volume 46. John Wiley & Sons, 2004.
- [6] Aapo Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Tr. Neur. Netw.*, 10(3):626–634, 1999.
- [7] Jerome H Friedman. Exploratory projection pursuit. *J. ASA*, 82(397):249–266, 1987.
- [8] Tijl De Bie. An information theoretic framework for data mining. In *Proc. of KDD*, pages 564–572. ACM, 2011.
- [9] Tijl De Bie. Subjective interestingness in exploratory data mining. In *Proc. of IDA*, pages 19–31. Springer, 2013.
- [10] Martin J Wainwright and Michael I Jordan. Graphical models, exponential families, and variational inference. *Found. Trends in ML*, 1(1-2):1–305, 2008.
- [11] Konstantinos Zografos. On maximum entropy characterization of pearson’s type II and VII multivariate distributions. *J. Multiv. Anal.*, 71(1):67–75, 1999.
- [12] Samuel Kotz and Saralees Nadarajah. *Multivariate t distributions and their applications*. Cambridge University Press, 2004.
- [13] Michael Roth. On the multivariate t distribution. Technical Report LiTH-ISY-R-3059, Department of Electrical Engineering, Linköping universitet, 2013.
- [14] Thomas M Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE Tr. Electr. Comp.*, pages 326–334, 1965.
- [15] Gene H Golub and Charles F Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.
- [16] Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *J. ACM*, 58(3), 2011.
- [17] Jiashi Feng, Huan Xu, and Shuicheng Yan. Robust PCA in high-dimension: A deterministic approach. *CoRR*, abs/1206.4628, 2012.